# Predicting Election Trends By Web Data
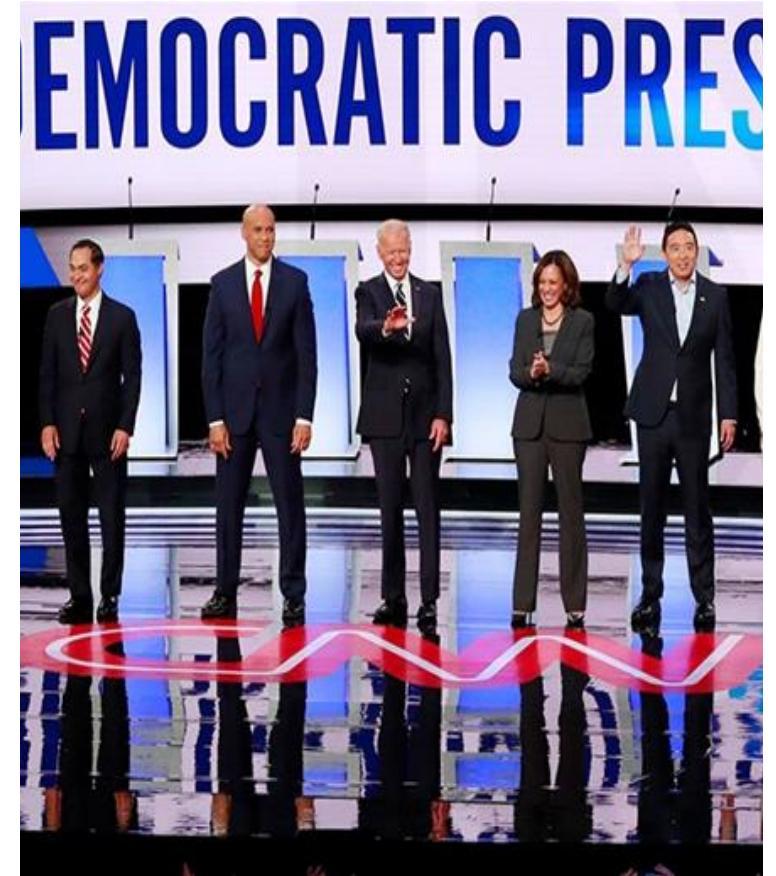
Kaan Cem Ketenci

Deji Suolang

# Research Questions

Can we predict dynamic changes in voter support of political candidates by social media and web data?

- How does predictions informed by real-time web-data map with existing benchmark national opinion polls?

- Is web data informative of shifts in public opinion before, during and after debate?

Focus: October 15, 2019  Democratic Debate

Larger volume of social media data/post/search volume/opinion shifts

# Our study will use four primary data sources:

1. **Twitter streaming API**

2. **2016 American Community Survey Data**

3. **Google Trends**

4. **Odd checkers Website**

# Twitter Trendsmap of Trending Topics During the Debate

https://www.trendsmap.com/

# Data source 1: Twitter streaming API

- Listening tweets from October 10 to October 20, 2019 for 10 days.

- Hashtags/Keywords: names and related terms of 12 candidates

- Tracking the number of twitter follower during the debate

# Twitter streaming API (cont'd)

**Sample tweets**

As we obtained a very large json file (100 GB in total), we randomly sampled 10000 for each of 12 candidate per day

**Segregating positive and negative tweets**

- Cleaning: Select relevant variables of the tweets such as "text", "retweet text" and "location", remove hashtags and URLs and other twitter handles.
- Use get_sentiment() function to extract sentiment score for each of the tweets to understand the change over time.
- Classify the tweet sentiment as positive, negative and neutral(1, -1 and 0).
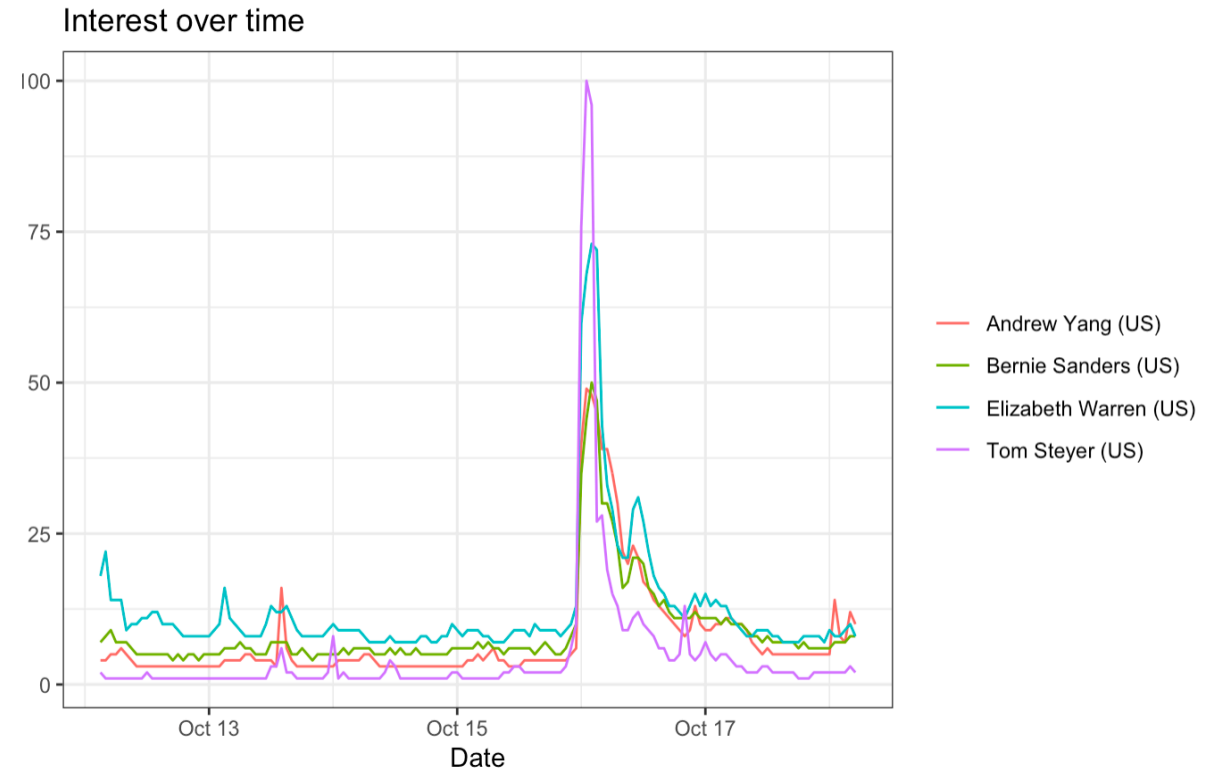
**Volumetric Information**

Changes in the number of followers and the public support of candidates may be correlated.

# Data Source 2: 2016 American Community Survey Data

- Use the "censusapi" R package to interface with the Census Bureau's API to streamline the data collection process

- Variable selections: Select socio-demographic information at state level from 2016 ACS including "age", "household income" and racial information to compute "Black African American Ratio".

- Data cleaning and recoding, for example, recode "66666" into NA

- Use "state name" variable to link with twitter data and google trends data for further analysis.

# Data Source 3: Google Trend

- Webscraping of daily and hourly number of Google searches of each of the 12 candidates' names between October 10 and October 20.

- This Google search statistic may primarily reflect dynamic shifts in popularity of each candidate and how often people are talking about them. It does not involve any sentiment analysis or categorization of searches into positve or negative searches for candidates.



- As expected, there is a big spike in the hourly searches in the evening of October 15th during the debate hours and several hours following the debate.

# Data Source 4: Odds Checker Website



https://www.oddschecker.com/politics/us-politics/us-presidential-election-2020/democrat-candidate

# Web Scraping from Odd checkers Website

- Oddschecker is a website that combines betting odds for various bets from all major betting companies. We are interested in the bets on who will become the Democratic Party US presidential nominee. We focus on 12 candidates who participated in the CNN debate and collected data between October 10 and October 20.

- Convert betting odds data into implied probabilities.

- We collected data once a day. We further collected data shortly before and after the debate.

- The increases and decreases in implied probabilities of winning the nomination can be compared with political analyses by mainstream news organizations of debate performances of each candidate.

# Questions?